ABSTRACT
        The purpose of this study was to investigate the effect of a
Rasch-based procedure to calibrate responses for funding applications. The
data set included 112 proposals and 66 readers, who independently scored
randomly assigned proposals using a scoring instrument. The data were
analyzed using FACETS (Linacre, 1999). The analysis indicated that the
instrument succeeded in separating proposals into four distinct strata of
quality. The proposal quality measures were found to be trustworthy in terms
of their accuracy and stability. The 17 items on the scoring instrument were
functioning as intended with a reliability of 0.98. Readers were found to be
internally consistent while using the instrument to assess the quality of the
proposals. Decision makers could easily translate the findings provided by
the analyses and rely on the precision of the estimates to make reliable and
defensible funding decisions. (Contains 5 tables and 11 references.) (Author)

ED 476 186

# A RASCH MEASUREMENT EXAMPLE IN GRANT APPLICATION PROCESS

Yesim CAPA

The Ohio State University


William E. LOADMAN

The Ohio State University

BEST COPY AVAILABLE

TM034910

Paper presented at the annual meeting of the American Educational Research Association,

Session 56.056: *Rasch Model Scaling*. April 24, 2003. Chicago, IL.

# Abstract

The purpose of this study was to investigate the effect of a Rasch-based procedure to calibrate responses for funding applications. The data set included 112 proposals and 66 readers, which independently scored randomly assigned proposals using a scoring instrument. The data were analyzed using *Facets* (Linacre, 1999).

The analysis indicated that the instrument succeeded in separating proposals into four distinct strata of quality. The proposal quality measures were found to be trustworthy in terms of their accuracy and stability. The 17 items on the scoring instrument were functioning as intended with a reliability of .98. Readers were found to be internally consistent while using the instrument to assess the quality of the proposals. Finally, decision makers could easily translate the findings provided by the analyses and rely on the precision of the estimates to make reliable and defensible funding decisions.

## A Rasch measurement example in grant application process

When complex assessment systems are used to allocate dollars for grants, based on a request for a proposal, decision makers need information that would help them to determine if all facets of the system of proposal review including the proposal reviewers (readers), the items, the proposals, and the rating scales are working as intended. It is important to simultaneously estimate the rated quality of the proposals, and the difficulty/ease of the items, and the severity/leniency of the readers. This could not be accomplished by a classical test theory approach. Utilizing the many-facet Rasch model allows for the simultaneous estimation of variations in the severity of judges, the item ease/difficulty, and the proposal quality. Thus, it eliminates or accounts for the differences in the measurement model beyond reader ratings of the reviewed proposal (Linacre, Wright, & Lunz, 1990).

Another advantage provided by the facet model is that its use does not necessitate each proposal to be reviewed by every reader on all of the items. Wright and Stone (1979) stated that the only requirement is a network to link every parameter to the other parameter by means of ratings. This network allows all measures to be placed on a common continuum.

The context for this study was a state initiative announced through a request for a proposal (RFP) initiative for competitively based funding support. The applicants (school districts) developed and submitted a written proposal seeking support for their proposal initiative. Multiple awards were anticipated in the announced RFP.

The purpose of this study was to investigate the effect of Rasch-based procedures to calibrate responses for funding applications. More specifically, it was aimed to estimate three sources of variability within the rating system: applications, items, and readers. The data collected in this study were analyzed by using Linacre's (1999) FACETS computer program.

4

In this paper, a brief theory related to the many-facet Rasch measurement will be presented to better inform the audience about the analyses that were performed. Then, statistical analyses and the terminology used by *Facets* will be explained in the context of present study.

The final section of the report will include subsections of the result part. First, the map will be examined because it allows us to see all facets of the analysis within a single frame of reference. Next, several questions will be addressed regarding the sources of variability: How well did applications calibrate? How well did readers calibrate? How well did items calibrate? Finally, a recommendation will be made to help decision makers use the information in the administration of the grant process.

*Many – Facet Rasch Model*

The many-facet analysis used in this study describes the probability that a specific application (n) rated by a specific reader (j) will be rated in a particular category (k) on a specific item (i). Here is the mathematical form of the three-facet Rasch model that shows the relationships among these facets in terms of a logistics odds ratio:

$$\text{Log}\left(\frac{P_{nijk}}{P_{nijk-1}}\right) = B_n - D_i - C_j - F_k$$

where

$P_{nijk}$ is the probability of application n being rated by reader j on item i with a rating of k,

$P_{nijk-1}$ is the probability of application n being rated by reader j on item i with a rating of k-1,

$B_n$ is the quality of the application n,

$D_i$ is the difficulty of item i,

$C_j$ is the severity of reader j,

$F_k$ is the difficulty (F) of category k of the rating scale (Linacre, 1994).

In this study, the measurement model specifies that a common rating scale category structure applies across all items and for all readers; in other words $F_k$ is constant across items and readers. It must be noted that there is no mathematical limit on the number of facets included in the model; however, most applications do not go beyond two or three facets in addition to the item and person (or application) facet (Smith, 1996).

The psychometric model, presented above in the mathematical form, includes three facets – applications, readers, and items. The *Facets* program uses the ratings that readers give on all items to estimate measures for each element of each facet. For each application, the measure is an estimate of that application's quality. The larger the measure, the better the application. For each reader, the measure is an estimate of the degree of severity that each reader exercised while evaluating the quality of the application. The larger the measure, the more severe the reader. For each item, the measure is an estimate of the difficulty of the item. The larger the measure, the more difficult it is for an application to obtain high ratings on that item. All of the measures that Facet produces are in the same linear unit of measure which is in logits, or log-odds units. Thus, the comparison within and between the facets of the analysis is straightforward and more importantly consistent (Myford & Wolfe, 2000).

Moreover, along with each measure the *Facets* program produces a standard error, which gives information about the precision of the logit estimate, and fit statistics, which provides information about how well the data fit the expectations of the measurement model. Fit statistics are presented by two measures: infit and outfit mean squares (Linacre, 1999b). Mean square, which is a chi-square statistic with an expectation of 1 and range of 0 to infinity, is based on the ratio of observed error variance to modeled error variance. The *Outfit* statistics are unweighted mean square residuals that are specifically sensitive to the outliers. On the other hand, the *Infit*

statistics weight each standardized residual by its variance and are more sensitive to unexpected responses near the point where decisions are made. A standardized value of the mean square statistics is also provided by the *Facets*. Different researchers have been using different cutoffs for identifying misfitting items, applications, or readers. Even some have been using standardized residuals. Wright, Linacre, Gustafson, and Martin (1994) reported that there are "no hard-and fast rules" for these measures and the decision depends on the purpose of research and the researcher (p.370). For example, high-stakes tests would tolerate less noise than low-stakes tests. For the purpose of this study, mean square values between 0.5 and 1.8 were investigated. The standardized values weren't used because they are affected by sample size.

It must be emphasized that in order to employ a Facets Model, the data must meet two requirements. First, the data must be approximately unidimensional, i.e., most of the items should produce data along the same underlying construct. Second, the data must show local independence, i.e., the probability of responding to one item should not affect the response to another item (Smith, 1996; Wright, 1996). These assumptions are common in all item response theory models.

## Method

The data set included 112 proposals (including a calibration application) and 66 readers. The calibration application, which was scored by all readers, was not an actual application but a proposal used to assist in the calibration of all readers. Readers underwent a training program, in which they are informed about the rating instrument and process, perspectives on reading and scoring applications, and data analysis. At least three readers were assigned randomly to evaluate the quality of the proposals. Readers independently scored each assigned proposal using the

scoring instrument. The trio of readers for each application was unique to each application as the three member composition of each trio was consistently rotated.

*Instrument*

A reader scoring instrument, generated on the basis of the content required in the RFP, consisted of 18 items of which 17 of them are used for final calibrated score. The 17 items include six-point rating scale: 1 is given when no evidence is provided; 4 is given when the item is addressed, and evidence is detailed with few examples of quality; and 6 is used when item is addressed, exceptionally well-developed and high-quality examples are presented. The last item assesses the overall quality of the application and includes six-point scale where 1 is the poor application and 6 is the exceptional application.

*Assumption check*

The unidimensionality assumption is likely to be satisfied because of the characteristics of the instrument. The correlation matrix among items (See Table 1) also indicates that all items are significantly correlated; this could be accepted as a sign that the instrument has a single dimension. Moreover, it should be noted: Wright suggested that violating this assumption would not cause significant difficulty. According to Hambleton, Swaminathan, & Rogers (1991), this assumption cannot be strictly met because of several external factors. What is required for the unidimensionality assumption to be met adequately by a set of test data is the presence of a "dominant" factor.

--------------------------------------------------------------------------------

Insert Table 1

--------------------------------------------------------------------------------

In order to not to violate the assumption of local independence, the item 18, assessing the overall quality of the proposal, was excluded from the *Facets* analysis.

*Missing Data*

In *Facets* analysis, there is considerable allowance for missing data. Each parameter could be estimated from the subset of observations. Linacre (1994) stated that there is no need for imputation of values in order to handle missing data. Analysis will produce stable estimates unless there is tremendous amount of missing data or all observations are in the same extreme category.

When inspecting data by application, it appears that there are only three cases in which three out of 17 responses are missing. This corresponds to approximately 18% of the data for that individual application, which could be ignored as a problem.

*Data Analysis*

The data were analyzed by using the FACETS program (Linacre, 1999a), a Rasch-based computer program based on an extension of Wright and Masters' rating scale model. The output from the FACETS analysis provided information about the calibrated quality of each proposal, the utility of each item, and the scoring behavior of each reader.

<div align="center">Results</div>

Before interpreting the results, the subset connection should be controlled. In this analysis, the subset connection was satisfied, which indicates that the estimated parameters of three facets can be placed on a common scale with the same origin and that there is no identification problem in the estimation process.

In this study, there are three facets: 112 applications, 66 readers, and 17 evaluation items. There were at least three readers per application. The Facets can accommodate unequal number

of readers per application. After producing the iteration report (4 PROX and 24 UCON

Iterations), it provided the list of unexpected responses – 32 unexpected residuals including 6

from the calibration application, which forms only .004% of the total responses (Table 2).

---------------------------------------------------------------------------

Insert Table 2

---------------------------------------------------------------------------

The map shown as Figure 1 displays all facets of the analysis in one figure; for that

reason, it is powerful to examine this map first before looking at the fits statistics for the facets

(Myford & Wolfe, 2000). The scale of the map is in "logits" for all facets, which forms an equal-

interval scale and enables comparison of the results.

---------------------------------------------------------------------------

Insert Figure 1

---------------------------------------------------------------------------

The **first column** in the map depicts the logit scale. The **second column** indicates the

estimated quality of applications. Higher quality applications appear at the top of the column,

while lower quality applications appear at the lower end of the column. Each star represents two

applications, and a dot represents one application. These measures of estimate show a fairly

negatively skewed distribution with three very low quality applications.

The **third column** displays the level of leniency – severity of the readers. In this column,

more severe readers appear at the higher end of the column, and more lenient readers appear at

the lower end of the column. The reader severity measures show a symmetric distribution with a

narrower range than the application measures. The reader severity measures were calibrated

around the mean (0), with 32 out of 66 readers on the severe side (logit higher than 0) and 28 out

of 66 readers on the lenient side (logit less than 0).

The **fourth column** shows the 17 items in terms of their relative ease – difficulty. Items appearing higher in the column are more difficult for applications to receive high ratings than the items appearing lower in the column. Item 17 appears to be the most difficult item in the instrument; the other items are relatively easier items.

From this point, the paper addresses itself to the detailed description of the results of application analysis, the reader analysis, and the item analysis.

*Application analysis*

Application measures are presented in ascending order of quality in Table 3.

-----------------------------------------------------------------------------------

Insert Table 3

-----------------------------------------------------------------------------------

Applications #96, #22, and #19 are the higher quality applications, whereas application #61 and #69 are the lower quality applications. The following information is presented in this table: observed score, observed count (the number of ratings), observed average (average for the ratings), and the fair average (average adjusted for reader severity), the logit measure and standard error, and the fit information including the infit and outfit mean square statistics for each application. By using the criteria previously established (logit score higher than 1.8 or lower than 0.5), when both infit and outfit statistics are examined, none of the applications present an area of concern. At the bottom of the table, overall statistics are presented. For this run, the separation reliability is .94, indicating that the differences among application measures are mainly due to the actual differences rather than the measurement error. The applications can also be separated into 4 distinct strata of quality (separation index = 4.11). The Root Mean

Square Error (RMSE) of 0.11 indicates a relatively low error in application measures.

Unlike the standard error of measurement in classical test theory, *Facets* provides a

separate estimate of standard error for each application. The standard error of measurement

indicates how much we would expect an application's quality estimate to change if different

readers and/or different items were used. The average standard error of measurement for the

applications is 0.11.

*Reader analysis*

Reader measures of severity/leniency are presented in Table 4.

-------------------------------------------------------------------------------

Insert Table 4

-------------------------------------------------------------------------------

Readers #1061 and #1013 are relatively more severe, and the readers #1038 and #1026

are relatively more lenient. By using the previously established criteria when both infit and outfit

statistics are examined, 3 out of 66 readers had either high (>1.8) or low (<0.5) infit and outfit

statistics. The readers that have low infit/outfit statistics are: #1031 and #1041. Their ratings tend

to be "muted." The low scores are of less concern, which are tending to show a flat line pattern

and little variation. In some cases, one might question whether the readers rate each item

independently or whether a halo effect or centrality effect may be operating (Myford &Wolfe,

2000).

The reader with high infit/outfit statistics is: #1066. When we check the unexpected

response table (Table 2), it appears that this reader rated the quality of application #22

unexpectedly low on items #8, #9, #11, and #13, while the application's rating is higher on

average. Moreover, this reader tended to rate leniently overall, therefore unexpectedly low rating

is surprising. This reader could be removed from the calibration of the application or provided with additional training for future use as a reader. The ratings of the readers with high infit measures tend to be "noisy," indicating that their rating shows more variation than expected in their ratings. Generally, it is recommended that readers with high infit/outfit statistics be trained for better quality of rating (Myford &Wolfe, 2000).

At the bottom of the table, overall statistics are presented. For this run, the separation reliability for readers is .91, indicating that the analysis is fairly reliable in separating readers into three different levels of severity and leniency (separation = 3.10). The RMSE score of .09 indicates that reader error is fairly low. To investigate whether the readers differ in their severity with which they rate applications, the fixed chi-square test is used. The chi-square of 874.6 with 65 degrees of freedom is significant (alpha=.01). This implies that readers are not considered equally severe/lenient after allowing for measurement error. This degree of leniency – severity is used to adjust the calibrated score for each application.

*Item Analysis*

Item measures are presented in Table 5. Items #17 and #15 are difficult to endorse, and items #3, #6, and #14 are relatively easier items in the instrument. These findings are also presented in Figure 1.

---------------------------------------------------------------------------------

Insert Table 5

---------------------------------------------------------------------------------

The *Facets* also provides a number of indications of the magnitude of the differences among elements of this facet, which is the difficulty/easiness of the items. These are: RMSE, Reliability, Separation Index, Fixed and Random Chi-square, Infit and Outfit statistics.

The RMSE score of 0.04 indicates that item error is very low. The reliability of items is very high (.98), and this is preferable for most studies. We could also interpret reliability such that the analysis is reliably separating items into approximately seven levels of difficulty (separation=6.68). Fixed chi-square tests the null hypothesis that all of the elements of item facet are equal. The fixed chi-square of 819.2 with degrees of freedom 16 is significant at the alpha .01, showing that the null hypothesis is rejected. Therefore, it could be concluded that items are not of equal difficulty/easiness. Random chi-square tests the hypothesis: "Can this set of elements be regarded as a random sample from a normal distribution?" The random chi-square of 16.0 with degrees of freedom 15 is not significant at the alpha .05. Thus, it could be implied that items could be regarded as a random sample from a normal distribution.

Using the same criteria to investigate the infit and outfit statistics (lower than 0.5 indicates muting; higher than 1.8 indicates noise), it appears that infit/outfit mean square of item #15 and #17 are higher than 1.8. That shows these items were not compatible with the quality estimates of the applications, and scores for these items may not be stable. The infit statistics of item #16 is also high (1.7), although not higher than the established criteria. This might occur because all of them appear as the last items. The readers might not have enough time to rate them adequately. Moreover, when we check the unexpected response table (Table 2), 14 out of 32 unexpected responses include item #15, #16, and particularly #17. The next time, the ordering of the questions might be changed to see if it happens again or the item might be modified for clarity. Because item #17 also appears as the most difficult item to endorse, this item could be revised or removed from the instrument.

## Conclusion

By utilizing the *Facets* analysis, we obtained specific information about how each element of each facet (i.e., each application, reader, and item) was performing. Results from our study indicated that the test succeeded in separating proposals (with a reliability of .94) into four distinct strata of quality. The distribution of the proposal quality measures was very similar in range to the distribution of reader severity measures. The *Facets* reports overall measure of accuracy and stability of proposal quality measures that is similar to the concept of standard error of measurement in classical test theory. The average standard error of measurement for the proposals was 0.11, which indicates fairly stable estimates of proposal quality if different readers or items were used. None of the proposals had infit and outfit mean-square indices either lower than 0.5 or higher than 1.8, which indicates the consistency shown in evaluating the quality of applications across items and across readers.

For items, the reliability was very high (.98), which is preferable in research studies. The infit and outfit mean-square indices for 17 items ranged from 0.6 to 2.1. Two of the items had fit statistics higher than 1.8, indicating noise or excess variation. Both of these items appear as the last items. Overall, it could be implied that rating on the items could be meaningfully combined to produce a single composite score to reflect the quality of the application.

The *Facets* also yielded a measure of the degree of severity each reader exercised while evaluating the proposals. The reader severity measures ranged from -0.66 logits to 0.66 logits, a 1.32 logit spread. The resulting chi-square value for readers was 945.1 with 65 degrees of freedom. This implied that readers could not be considered as equally severe/lenient after allowing for the measurement error. In addition, reader fit statistics provided evidence that

readers were internally consistent while using the instrument to assess the quality of the application.

In general, it would be concluded that all of the elements of the study including applications, readers, and items were functioning as intended. The Rasch measurement has an effective role to play in analysis and reporting of the data collected for the purpose of grant application review. The decision makers could easily use the findings provided by the analysis and rely on the accuracy of the estimates considering that the severity of the readers are also taken into account in the measurement model. Moreover, implications could be drawn for better rating process.

References

Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage Publications.

Linacre, J. M. (1999a). *Facets*, Version 3.22 [Computer program]. Chicago: MESA Press.

Linacre, J. M. (1999b). *A user's guide to Facets: Rasch measurement computer program*. Chicago: MESA Press.

Linacre, J. M. (1994). *Many-facet Rasch measurement*. Chicago: MESA Press.

Linacre, J. M., Wright, B. D., & Lunz, M. E. (1990). *A facets model for judgmental scoring*. (MESA Memo No. 61). Chicago: MESA Press.

Myford, C. M., & Wolfe, E. W. (2000). *Strengthening the ties that bind: Improving the linking network in sparsely connected rating designs*. (TOEFL Technical Report No.15). Princeton, NJ: Educational Testing Service.

Smith, R. M. (1996). A comparison of methods for dimensionality in Rasch measurement. *Structural Equation Modeling, 3*, 25-40.

Wright, B.D. (1996). Comparing Rasch Measurement and Factor Analysis. *Structural Equation Modeling, 3*, 3 – 24.

Wright, B. D., Linacre, J. M., Gustafson, J.E., & Martin-Löf P. (1994). Reasonable mean-square fit values. *Rasch Measurement Transactions, 8*(3), 370.

Wright, B.D., & Linacre J.M. (1985). *Microscale Manual*. Westport, Conn.: Mediax Interactive Technologies, Inc.

Wright, B. D., & Stone, M. H. (1979). *Best test design*. Chicago: MESA Press, University of Chicago.

Table 1

*Intercorrelations for Items*

| | I2 | I3 | I4 | I5 | I6 | I7 | I8 | I9 | I10 | I11 | I12 | I13 | I14 | I15 | I16 | I17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| I1 | 0.48 | 0.45 | 0.46 | 0.37 | 0.39 | 0.47 | 0.33 | 0.37 | 0.35 | 0.26 | 0.31 | 0.46 | 0.33 | 0.28 | 0.30 | 0.15 |
| I2 | | 0.51 | 0.50 | 0.58 | 0.50 | 0.60 | 0.41 | 0.39 | 0.38 | 0.35 | 0.29 | 0.36 | 0.33 | 0.22 | 0.29 | 0.18 |
| I3 | | | 0.67 | 0.57 | 0.54 | 0.62 | 0.39 | 0.43 | 0.44 | 0.41 | 0.37 | 0.46 | 0.36 | 0.23 | 0.33 | 0.26 |
| I4 | | | | 0.50 | 0.52 | 0.52 | 0.38 | 0.40 | 0.38 | 0.31 | 0.38 | 0.41 | 0.32 | 0.31 | 0.35 | 0.25 |
| I5 | | | | | 0.60 | 0.57 | 0.37 | 0.43 | 0.44 | 0.42 | 0.38 | 0.38 | 0.34 | 0.12 | 0.22 | 0.21 |
| I6 | | | | | | 0.58 | 0.40 | 0.44 | 0.44 | 0.40 | 0.36 | 0.37 | 0.34 | 0.26 | 0.28 | 0.20 |
| I7 | | | | | | | 0.48 | 0.47 | 0.45 | 0.42 | 0.43 | 0.39 | 0.42 | 0.25 | 0.34 | 0.21 |
| I8 | | | | | | | | 0.34 | 0.41 | 0.42 | 0.36 | 0.36 | 0.34 | 0.30 | 0.29 | 0.17 |
| I9 | | | | | | | | | 0.61 | 0.47 | 0.52 | 0.51 | 0.38 | 0.23 | 0.35 | 0.19 |
| I10 | | | | | | | | | | 0.61 | 0.52 | 0.43 | 0.44 | 0.13 | 0.28 | 0.18 |
| I11 | | | | | | | | | | | 0.53 | 0.44 | 0.39 | 0.18 | 0.28 | 0.22 |
| I12 | | | | | | | | | | | | 0.50 | 0.44 | 0.24 | 0.40 | 0.27 |
| I13 | | | | | | | | | | | | | 0.56 | 0.28 | 0.36 | 0.22 |
| I14 | | | | | | | | | | | | | | 0.24 | 0.27 | 0.24 |
| I15 | | | | | | | | | | | | | | | 0.54 | 0.25 |
| I16 | | | | | | | | | | | | | | | | 0.29 |

*Note.* All correlations are significant at the .05 significance level

18

Table 2

*Unexpected Responses (32 residuals sorted by order in data).*

| Cat | Step | Exp. | Resd | StRes | Num | app | Num | read | Nu | it |
|-----|------|------|------|-------|-----|-----|-----|------|----|----|
| 2 | 2 | 5.0 | -3.0 | -3 | 11 | 11 | 1039 | 1039 | 1 | 1 |
| 1 | 1 | 4.7 | -3.7 | -3 | 11 | 11 | 1039 | 1039 | 15 | 15 |
| 1 | 1 | 4.9 | -3.9 | -3 | 11 | 11 | 1039 | 1039 | 16 | 16 |
| 3 | 3 | 5.5 | -2.5 | -3 | 19 | 19 | 1064 | 1064 | 2 | 2 |
| 3 | 3 | 5.5 | -2.5 | -3 | 19 | 19 | 1064 | 1064 | 5 | 5 |
| 2 | 2 | 5.5 | -3.5 | -4 | 22 | 22 | 1066 | 1066 | 8 | 8 |
| 3 | 3 | 5.5 | -2.5 | -3 | 22 | 22 | 1066 | 1066 | 9 | 9 |
| 3 | 3 | 5.4 | -2.4 | -3 | 22 | 22 | 1066 | 1066 | 11 | 11 |
| 3 | 3 | 5.5 | -2.5 | -3 | 22 | 22 | 1066 | 1066 | 13 | 13 |
| 1 | 1 | 4.8 | -3.8 | -3 | 25 | 25 | 1016 | 1016 | 3 | 3 |
| 1 | 1 | 4.7 | -3.7 | -3 | 34 | 34 | 1047 | 1047 | 13 | 13 |
| 2 | 2 | 5.1 | -3.1 | -3 | 47 | 47 | 1024 | 1024 | 14 | 14 |
| 1 | 1 | 4.6 | -3.6 | -3 | 48 | 48 | 1013 | 1013 | 12 | 12 |
| 1 | 1 | 4.7 | -3.7 | -3 | 49 | 49 | 1059 | 1059 | 16 | 16 |
| 2 | 2 | 5.1 | -3.1 | -3 | 53 | 53 | 1002 | 1002 | 12 | 12 |
| 6 | 6 | 2.2 | 3.8 | 3 | 55 | 55 | 1011 | 1011 | 17 | 17 |
| 4 | 4 | 1.5 | 2.5 | 3 | 69 | 69 | 1021 | 1021 | 12 | 12 |
| 1 | 1 | 4.6 | -3.6 | -3 | 87 | 87 | 1059 | 1059 | 16 | 16 |
| 6 | 6 | 2.1 | 3.9 | 3 | 88 | 88 | 1052 | 1052 | 17 | 17 |
| 6 | 6 | 2.1 | 3.9 | 3 | 92 | 92 | 1016 | 1016 | 17 | 17 |
| 6 | 6 | 2.0 | 4.0 | 3 | 94 | 94 | 1058 | 1058 | 17 | 17 |
| 1 | 1 | 4.6 | -3.6 | -3 | 96 | 96 | 1013 | 1013 | 8 | 8 |
| 1 | 1 | 5.1 | -4.1 | -4 | 106 | 106 | 1038 | 1038 | 17 | 17 |
| 5 | 5 | 1.8 | 3.2 | 3 | 108 | 108 | 1040 | 1040 | 15 | 15 |
| 4 | 4 | 1.5 | 2.5 | 3 | 108 | 108 | 1040 | 1040 | 17 | 17 |
| 1 | 1 | 4.6 | -3.6 | -3 | 110 | 110 | 1068 | 1068 | 16 | 16 |
| 2 | 2 | 5.0 | -3.0 | -3 | 5000 | 5000 | 1033 | 1033 | 5 | 5 |
| 1 | 1 | 4.8 | -3.8 | -3 | 5000 | 5000 | 1038 | 1038 | 15 | 15 |
| 1 | 1 | 4.9 | -3.9 | -3 | 5000 | 5000 | 1048 | 1048 | 14 | 14 |
| 2 | 2 | 5.1 | -3.1 | -3 | 5000 | 5000 | 1064 | 1064 | 14 | 14 |
| 1 | 1 | 4.6 | -3.6 | -3 | 5000 | 5000 | 1069 | 1069 | 16 | 16 |
| 2 | 2 | 5.1 | -3.1 | -3 | 5000 | 5000 | 1008 | 1008 | 14 | 14 |
| Cat | Step | Exp. | Resd | StRes | Num | app | Num | read | Nu | it |

```
Vertical = (1*,2*,3A) Yardstick (columns,lines,low,high)= 0,10,-2,2
------------------------------------------------------------
|Measr|+app     |-readers|-items                  |S.1 |
------------------------------------------------------------
+   2 +          +        +                        +(6)  +
|     |          |        |                        |     |
|     | High     | Severe | Difficult              |     |
|     | score    |        |                        |     |
|     |          |        |                        | --- |
|     |          |        |                        |     |
|     |          |        |                        |     |
|     |          |        |                        |     |
|     |          |        |                        |     |
|     | **       |        |                        |     |
+   1 + **       +        +                        +     +
|     | *.       |        |                        | 5   |
|     | ***.     |        |                        |     |
|     | **       |        | 17                     |     |
|     | ****.    | .      |                        |     |
|     | *****    | *.     |                        | --- |
|     | *****.   | **     | 15                     |     |
|     | ****     | *      |                        |     |
|     | ****     | *****  | 11   16                | 4   |
|     | *****.   | ****** | 8    10                |     |
*   0 *  *****.  * ***    * 1                      *     *
|     | ***      | ****.  | 2    5    9   12   13  | --- |
|     | **       | *.     | 4    7                 |     |
|     | *        | ***.   | 3    6    14           | 3   |
|     | *.       | ***    |                        |     |
|     | *        | *      |                        | --- |
|     | *        | .      |                        |     |
|     | .        |        |                        |     |
|     |          |        |                        | 2   |
+  -1 +          +        +                        +     +
|     | .        |        |                        |     |
|     |          |        |                        |     |
|     |          |        |                        |     |
|     | .        |        |                        | --- |
|     |          |        |                        |     |
|     |          |        |                        |     |
|     | Low      | Lenient| Easy                   |     |
|     | score    |        |                        |     |
+  -2 +          +        +                        +(1)  +
------------------------------------------------------------
|Measr| * = 2   | * = 2  |-items                   |S.1 |
------------------------------------------------------------
```

Figure 1

*Variable Map of Application Measures, Reader Measures, and Item Calibrations*

Table 3

*Summary Statistics for Application Measures Ranked by the Logit Scores*

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | ZStd | Outfit MnSq | ZStd | Num | app |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 85 | 51 | 1.7 | 1.49 | −1.46 | .15 | 0.6 | −1 | 0.6 | −1 | 61 | 61 |
| 95 | 51 | 1.9 | 1.79 | −1.10 | .14 | 1.1 | 0 | 1.1 | 0 | 69 | 69 |
| 111 | 51 | 2.2 | 2.34 | −.68 | .12 | 1.1 | 0 | 1.3 | 1 | 108 | 108 |
| 158 | 50 | 3.2 | 2.43 | −.63 | .10 | 1.2 | 1 | 1.3 | 1 | 94 | 94 |
| 174 | 68 | 2.6 | 2.46 | −.61 | .09 | 0.6 | −3 | 0.6 | −3 | 6 | 6 |
| 190 | 68 | 2.8 | 2.63 | −.52 | .09 | 0.7 | −1 | 0.7 | −2 | 5 | 5 |
| 151 | 51 | 3.0 | 2.74 | −.46 | .10 | 1.1 | 0 | 1.2 | 0 | 64 | 64 |
| 134 | 51 | 2.6 | 2.77 | −.44 | .11 | 0.6 | −2 | 0.6 | −2 | 71 | 71 |
| 156 | 51 | 3.1 | 2.86 | −.40 | .10 | 0.8 | −1 | 0.8 | −1 | 109 | 109 |
| 161 | 51 | 3.2 | 2.90 | −.38 | .10 | 0.6 | −2 | 0.6 | −2 | 105 | 105 |
| 137 | 51 | 2.7 | 3.02 | −.32 | .11 | 0.6 | −2 | 0.6 | −2 | 44 | 44 |
| 160 | 51 | 3.1 | 3.05 | −.30 | .10 | 1.2 | 1 | 1.2 | 1 | 73 | 73 |
| 237 | 68 | 3.5 | 3.19 | −.23 | .09 | 0.6 | −3 | 0.6 | −3 | 27 | 27 |
| 193 | 51 | 3.8 | 3.27 | −.19 | .10 | 1.3 | 1 | 1.2 | 1 | 40 | 40 |
| 197 | 49 | 4.0 | 3.33 | −.17 | .11 | 0.6 | −2 | 0.6 | −2 | 111 | 111 |
| 159 | 50 | 3.2 | 3.35 | −.16 | .10 | 0.6 | −3 | 0.6 | −3 | 60 | 60 |
| 187 | 51 | 3.7 | 3.37 | −.15 | .10 | 0.9 | 0 | 0.9 | 0 | 74 | 74 |
| 177 | 51 | 3.5 | 3.43 | −.12 | .10 | 0.8 | −1 | 0.8 | −1 | 78 | 78 |
| 252 | 68 | 3.7 | 3.47 | −.10 | .09 | 0.8 | −1 | 0.8 | −1 | 4 | 4 |
| 190 | 51 | 3.7 | 3.48 | −.09 | .10 | 0.9 | 0 | 0.9 | 0 | 101 | 101 |
| 199 | 51 | 3.9 | 3.51 | −.08 | .11 | 1.0 | 0 | 1.0 | 0 | 56 | 56 |
| 195 | 51 | 3.8 | 3.56 | −.05 | .10 | 1.0 | 0 | 1.0 | 0 | 77 | 77 |
| 248 | 68 | 3.6 | 3.68 | .01 | .09 | 1.0 | 0 | 0.9 | 0 | 41 | 41 |
| 196 | 51 | 3.8 | 3.71 | .02 | .11 | 0.7 | −1 | 0.7 | −2 | 42 | 42 |
| 206 | 51 | 4.0 | 3.70 | .02 | .11 | 1.6 | 2 | 1.5 | 2 | 62 | 62 |
| 183 | 50 | 3.7 | 3.70 | .02 | .10 | 1.0 | 0 | 1.0 | 0 | 88 | 88 |
| 253 | 68 | 3.7 | 3.74 | .04 | .09 | 1.3 | 1 | 1.3 | 2 | 7 | 7 |
| 238 | 68 | 3.5 | 3.75 | .04 | .09 | 0.8 | −1 | 0.8 | −1 | 9 | 9 |
| 203 | 51 | 4.0 | 3.75 | .04 | .11 | 0.9 | 0 | 0.8 | −1 | 76 | 76 |
| 201 | 50 | 4.0 | 3.75 | .04 | .11 | 0.9 | 0 | 0.9 | 0 | 85 | 85 |
| 187 | 51 | 3.7 | 3.74 | .04 | .10 | 1.1 | 0 | 1.2 | 1 | 89 | 89 |
| 205 | 51 | 4.0 | 3.75 | .04 | .11 | 0.8 | −1 | 0.7 | −1 | 95 | 95 |
| 168 | 51 | 3.3 | 3.77 | .05 | .10 | 0.8 | −1 | 0.8 | −1 | 102 | 102 |
| 163 | 51 | 3.2 | 3.76 | .05 | .10 | 0.5 | −4 | 0.5 | −4 | 107 | 107 |
| 183 | 51 | 3.6 | 3.78 | .06 | .10 | 0.9 | 0 | 0.9 | 0 | 58 | 58 |
| 202 | 49 | 4.1 | 3.79 | .06 | .11 | 1.2 | 0 | 1.1 | 0 | 75 | 75 |
| 270 | 68 | 4.0 | 3.81 | .07 | .09 | 1.0 | 0 | 0.9 | 0 | 18 | 18 |
| 234 | 68 | 3.4 | 3.80 | .07 | .09 | 0.9 | 0 | 0.8 | −1 | 37 | 37 |
| 183 | 51 | 3.6 | 3.81 | .07 | .10 | 0.9 | 0 | 0.9 | 0 | 66 | 66 |
| 211 | 51 | 4.1 | 3.80 | .07 | .11 | 1.5 | 2 | 1.5 | 2 | 100 | 100 |
| 246 | 68 | 3.6 | 3.82 | .08 | .09 | 1.0 | 0 | 1.0 | 0 | 33 | 33 |
| 203 | 51 | 4.0 | 3.83 | .08 | .11 | 1.0 | 0 | 0.9 | 0 | 72 | 72 |
| 209 | 51 | 4.1 | 3.85 | .09 | .11 | 1.7 | 3 | 1.6 | 2 | 55 | 55 |
| 184 | 51 | 3.6 | 3.96 | .15 | .10 | 1.1 | 0 | 1.1 | 0 | 92 | 92 |
| 194 | 51 | 3.8 | 3.95 | .15 | .11 | 0.8 | −1 | 0.7 | −1 | 93 | 93 |
| 192 | 51 | 3.8 | 4.00 | .17 | .10 | 1.3 | 1 | 1.3 | 1 | 99 | 99 |
| 245 | 68 | 3.6 | 4.05 | .20 | .09 | 0.7 | −2 | 0.7 | −2 | 13 | 13 |
| 216 | 51 | 4.2 | 4.05 | .20 | .11 | 1.4 | 1 | 1.3 | 1 | 84 | 84 |
| 280 | 67 | 4.2 | 4.06 | .21 | .10 | 1.0 | 0 | 0.9 | 0 | 35 | 35 |
| 233 | 65 | 3.6 | 4.09 | .22 | .09 | 1.0 | 0 | 1.0 | 0 | 54 | 54 |
| 215 | 51 | 4.2 | 4.09 | .23 | .11 | 0.9 | 0 | 0.9 | 0 | 98 | 98 |
| 294 | 67 | 4.4 | 4.13 | .25 | .10 | 1.6 | 3 | 1.5 | 2 | 34 | 34 |
| 270 | 67 | 4.0 | 4.16 | .26 | .09 | 1.3 | 1 | 1.3 | 1 | 23 | 23 |
| 272 | 68 | 4.0 | 4.15 | .26 | .09 | 0.9 | 0 | 0.9 | 0 | 31 | 31 |

Table 3 (continued)

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | ZStd | Outfit MnSq | ZStd | Num | app |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 336 | 84 | 4.0 | 4.16 | .26 | .09 | 0.9 | 0 | 0.9 | 0 | 46 | 46 |
| 209 | 51 | 4.1 | 4.20 | .29 | .11 | 0.8 | -1 | 0.8 | -1 | 70 | 70 |
| 217 | 51 | 4.3 | 4.22 | .30 | .11 | 0.7 | -1 | 0.6 | -2 | 90 | 90 |
| 233 | 51 | 4.6 | 4.27 | .33 | .12 | 0.9 | 0 | 0.9 | 0 | 65 | 65 |
| 230 | 51 | 4.5 | 4.27 | .33 | .12 | 0.9 | 0 | 0.8 | -1 | 81 | 81 |
| 294 | 68 | 4.3 | 4.29 | .34 | .10 | 0.9 | 0 | 0.9 | 0 | 1 | 1 |
| 297 | 67 | 4.4 | 4.31 | .35 | .10 | 1.2 | 0 | 1.1 | 0 | 36 | 36 |
| 200 | 51 | 3.9 | 4.33 | .36 | .11 | 1.2 | 0 | 1.1 | 0 | 59 | 59 |
| 268 | 67 | 4.0 | 4.34 | .37 | .09 | 0.8 | -1 | 0.8 | -1 | 12 | 12 |
| 273 | 68 | 4.0 | 4.34 | .37 | .09 | 0.8 | -1 | 0.8 | -1 | 15 | 15 |
| 383 | 85 | 4.5 | 4.34 | .37 | .09 | 1.0 | 0 | 1.0 | 0 | 50 | 50 |
| 281 | 68 | 4.1 | 4.37 | .39 | .10 | 1.2 | 1 | 1.2 | 1 | 38 | 38 |
| **4654** | **1081** | **4.3** | **4.38** | **.40** | **.02** | **1.0** | **0** | **1.0** | **0** | **5000** | **5000** |
| 280 | 68 | 4.1 | 4.40 | .41 | .09 | 0.9 | 0 | 0.8 | 0 | 3 | 3 |
| 305 | 68 | 4.5 | 4.40 | .41 | .10 | 0.8 | -1 | 0.7 | -1 | 52 | 52 |
| 265 | 68 | 3.9 | 4.40 | .41 | .09 | 0.9 | 0 | 0.9 | 0 | 79 | 79 |
| 221 | 51 | 4.3 | 4.43 | .43 | .12 | 0.9 | 0 | 0.9 | 0 | 91 | 91 |
| 294 | 68 | 4.3 | 4.49 | .47 | .10 | 1.4 | 1 | 1.2 | 1 | 17 | 17 |
| 276 | 68 | 4.1 | 4.53 | .49 | .09 | 0.9 | 0 | 0.9 | 0 | 29 | 29 |
| 281 | 68 | 4.1 | 4.54 | .50 | .10 | 1.0 | 0 | 0.9 | 0 | 20 | 20 |
| 222 | 49 | 4.5 | 4.53 | .50 | .12 | 1.2 | 1 | 1.1 | 0 | 51 | 51 |
| 219 | 51 | 4.3 | 4.55 | .51 | .11 | 0.7 | -1 | 0.6 | -2 | 103 | 103 |
| 303 | 68 | 4.5 | 4.57 | .52 | .10 | 0.9 | 0 | 0.8 | 0 | 82 | 82 |
| 242 | 51 | 4.7 | 4.57 | .53 | .13 | 1.2 | 0 | 1.0 | 0 | 43 | 43 |
| 287 | 68 | 4.2 | 4.59 | .54 | .10 | 1.1 | 0 | 1.1 | 0 | 14 | 14 |
| 223 | 51 | 4.4 | 4.59 | .54 | .11 | 1.0 | 0 | 0.9 | 0 | 16 | 16 |
| 378 | 85 | 4.4 | 4.60 | .55 | .09 | 0.9 | 0 | 0.9 | 0 | 28 | 28 |
| 305 | 67 | 4.6 | 4.62 | .56 | .11 | 1.3 | 1 | 1.4 | 1 | 11 | 11 |
| 223 | 50 | 4.5 | 4.62 | .56 | .12 | 1.0 | 0 | 0.9 | 0 | 67 | 67 |
| 316 | 68 | 4.6 | 4.66 | .59 | .11 | 1.0 | 0 | 0.9 | 0 | 8 | 8 |
| 382 | 85 | 4.5 | 4.67 | .60 | .09 | 1.2 | 1 | 1.1 | 0 | 26 | 26 |
| 326 | 68 | 4.8 | 4.67 | .60 | .11 | 1.1 | 0 | 1.1 | 0 | 47 | 47 |
| 225 | 50 | 4.5 | 4.66 | .60 | .12 | 1.3 | 1 | 1.2 | 0 | 86 | 86 |
| 328 | 68 | 4.8 | 4.72 | .64 | .11 | 1.2 | 1 | 1.2 | 0 | 2 | 2 |
| 249 | 51 | 4.9 | 4.72 | .64 | .14 | 0.9 | 0 | 0.8 | -1 | 104 | 104 |
| 389 | 85 | 4.6 | 4.73 | .65 | .10 | 1.3 | 1 | 1.2 | 0 | 25 | 25 |
| 237 | 50 | 4.7 | 4.74 | .65 | .13 | 1.2 | 0 | 1.2 | 0 | 87 | 87 |
| 323 | 68 | 4.8 | 4.77 | .68 | .11 | 1.1 | 0 | 1.0 | 0 | 49 | 49 |
| 232 | 51 | 4.5 | 4.79 | .70 | .12 | 1.0 | 0 | 1.0 | 0 | 80 | 80 |
| 232 | 49 | 4.7 | 4.83 | .73 | .13 | 1.4 | 1 | 1.2 | 0 | 97 | 97 |
| 321 | 67 | 4.8 | 4.88 | .77 | .11 | 1.0 | 0 | 0.9 | 0 | 24 | 24 |
| 239 | 49 | 4.9 | 4.89 | .78 | .14 | 1.4 | 1 | 1.3 | 0 | 83 | 83 |
| 254 | 51 | 5.0 | 4.90 | .80 | .14 | 1.3 | 1 | 1.2 | 0 | 53 | 53 |
| 325 | 68 | 4.8 | 4.92 | .81 | .11 | 1.1 | 0 | 1.1 | 0 | 21 | 21 |
| 239 | 50 | 4.8 | 4.92 | .81 | .13 | 0.8 | 0 | 0.9 | 0 | 68 | 68 |
| 252 | 51 | 4.9 | 4.92 | .82 | .14 | 0.7 | -1 | 0.6 | -1 | 45 | 45 |
| 329 | 68 | 4.8 | 4.96 | .85 | .11 | 1.3 | 1 | 1.1 | 0 | 30 | 30 |
| 340 | 68 | 5.0 | 4.96 | .85 | .12 | 1.1 | 0 | 1.0 | 0 | 39 | 39 |
| 250 | 51 | 4.9 | 5.03 | .92 | .14 | 1.0 | 0 | 0.9 | 0 | 110 | 110 |
| 355 | 68 | 5.2 | 5.05 | .95 | .14 | 1.0 | 0 | 0.8 | 0 | 10 | 10 |
| 332 | 68 | 4.9 | 5.07 | .96 | .12 | 1.0 | 0 | 0.8 | -1 | 32 | 32 |
| 224 | 49 | 4.6 | 5.07 | .96 | .12 | 1.0 | 0 | 0.9 | 0 | 63 | 63 |
| 339 | 68 | 5.0 | 5.14 | 1.04 | .12 | 1.3 | 1 | 1.0 | 0 | 48 | 48 |
| 250 | 51 | 4.9 | 5.14 | 1.04 | .14 | 1.2 | 0 | 0.9 | 0 | 57 | 57 |
| 252 | 50 | 5.0 | 5.14 | 1.05 | .15 | 1.8 | 2 | 1.6 | 1 | 106 | 106 |
| 342 | 68 | 5.0 | 5.18 | 1.10 | .13 | 1.1 | 0 | 1.1 | 0 | 19 | 19 |
| 350 | 68 | 5.1 | 5.19 | 1.11 | .13 | 1.4 | 1 | 1.7 | 2 | 22 | 22 |

Table 3 (continued)

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | ZStd | Outfit MnSq | ZStd | Num | app |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 252 | 51 | 4.9 | 5.22 | 1.15 | .14 | 1.6 | 2 | 1.5 | 1 | 96 | 96 |
| 280.1 | 68.0 | 4.1 | 4.10 | .28 | .11 | 1.0 | -0.1 | 1.0 | -0.3 | Mean | |
| 419.9 | 96.7 | 0.7 | 0.75 | .46 | .02 | 0.3 | 1.5 | 0.2 | 1.4 | S.D. | |

RMSE (Model) .11 Adj S.D. .45 Separation 4.11 Reliability .94
Fixed (all same) chi-square: 1761.2 d.f.: 111 significance: .00
Random (normal) chi-square: 109.4 d.f.: 110 significance: .50

Table 4

*Summary Statistics for Reader Measures Ranked by the Number of Readers*

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | ZStd | Outfit MnSq | ZStd | Num | readers |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 416 | 102 | 4.1 | 4.38 | -.12 | .08 | 1.0 | 0 | 1.0 | 0 | 1002 | 1002 |
| 78 | 17 | 4.6 | 4.23 | -.03 | .21 | 0.6 | -1 | 0.6 | -1 | 1003 | 1003 |
| 635 | 136 | 4.7 | 4.65 | -.31 | .08 | 1.2 | 1 | 1.1 | 0 | 1006 | 1006 |
| 430 | 119 | 3.6 | 4.73 | -.37 | .07 | 1.0 | 0 | 1.0 | 0 | 1007 | 1007 |
| 382 | 85 | 4.5 | 4.71 | -.35 | .09 | 1.4 | 2 | 1.4 | 2 | 1008 | 1008 |
| 583 | 153 | 3.8 | 3.79 | .21 | .06 | 0.6 | -4 | 0.6 | -4 | 1009 | 1009 |
| 257 | 68 | 3.8 | 3.84 | .19 | .09 | 0.6 | -3 | 0.6 | -3 | 1010 | 1010 |
| 320 | 85 | 3.8 | 3.98 | .11 | .08 | 1.5 | 3 | 1.5 | 3 | 1011 | 1011 |
| 340 | 84 | 4.0 | 3.15 | .53 | .09 | 1.3 | 2 | 1.4 | 2 | 1013 | 1013 |
| 483 | 119 | 4.1 | 3.81 | .21 | .07 | 1.0 | 0 | 1.0 | 0 | 1014 | 1014 |
| 553 | 118 | 4.7 | 4.57 | -.25 | .08 | 1.0 | 0 | 0.9 | 0 | 1015 | 1015 |
| 457 | 119 | 3.8 | 3.74 | .24 | .07 | 1.2 | 2 | 1.3 | 2 | 1016 | 1016 |
| 523 | 115 | 4.5 | 4.37 | -.12 | .08 | 1.4 | 2 | 1.2 | 1 | 1017 | 1017 |
| 265 | 68 | 3.9 | 3.93 | .14 | .09 | 0.6 | -3 | 0.5 | -3 | 1018 | 1018 |
| 606 | 136 | 4.5 | 4.17 | .01 | .07 | 1.1 | 0 | 1.0 | 0 | 1019 | 1019 |
| 526 | 136 | 3.9 | 3.74 | .24 | .07 | 0.8 | -1 | 0.8 | -1 | 1020 | 1020 |
| 78 | 34 | 2.3 | 3.46 | .38 | .15 | 0.9 | 0 | 1.0 | 0 | 1021 | 1021 |
| 631 | 153 | 4.1 | 4.45 | -.17 | .06 | 0.8 | -2 | 0.8 | -1 | 1022 | 1022 |
| 432 | 100 | 4.3 | 4.41 | -.14 | .08 | 0.8 | -1 | 0.8 | -1 | 1024 | 1024 |
| 668 | 170 | 3.9 | 3.97 | .12 | .06 | 0.8 | -2 | 0.8 | -2 | 1025 | 1025 |
| 394 | 85 | 4.6 | 4.91 | -.53 | .10 | 0.6 | -2 | 0.6 | -2 | 1026 | 1026 |
| 390 | 102 | 3.8 | 3.87 | .17 | .07 | 0.6 | -4 | 0.6 | -3 | 1029 | 1029 |
| 592 | 153 | 3.9 | 4.02 | .09 | .06 | 0.7 | -2 | 0.7 | -3 | 1030 | 1030 |
| **339** | **102** | **3.3** | **3.44** | **.39** | **.08** | **0.5** | **-4** | **0.5** | **-4** | **1031** | **1031** |
| 623 | 136 | 4.6 | 4.63 | -.29 | .08 | 1.0 | 0 | 0.8 | -1 | 1032 | 1032 |
| 470 | 102 | 4.6 | 4.75 | -.39 | .09 | 1.4 | 2 | 1.4 | 2 | 1033 | 1033 |
| 492 | 114 | 4.3 | 4.09 | .05 | .08 | 1.6 | 3 | 1.5 | 3 | 1034 | 1034 |
| 530 | 118 | 4.5 | 4.47 | -.18 | .08 | 1.3 | 2 | 1.2 | 1 | 1035 | 1035 |
| 596 | 136 | 4.4 | 4.64 | -.30 | .07 | 0.8 | -1 | 0.8 | -1 | 1036 | 1036 |
| 616 | 152 | 4.1 | 4.22 | -.02 | .07 | 1.0 | 0 | 0.9 | 0 | 1037 | 1037 |
| 785 | 153 | 5.1 | 5.03 | -.64 | .09 | 1.2 | 1 | 1.0 | 0 | 1038 | 1038 |
| 638 | 135 | 4.7 | 4.79 | -.42 | .09 | 1.6 | 3 | 1.5 | 2 | 1039 | 1039 |
| 647 | 151 | 4.3 | 4.11 | .04 | .07 | 1.1 | 0 | 1.1 | 0 | 1040 | 1040 |
| **480** | **135** | **3.6** | **4.22** | **-.03** | **.06** | **0.5** | **-6** | **0.5** | **-6** | **1041** | **1041** |
| 374 | 85 | 4.4 | 4.04 | .08 | .09 | 1.0 | 0 | 0.9 | 0 | 1045 | 1045 |
| 364 | 85 | 4.3 | 4.35 | -.10 | .09 | 0.6 | -3 | 0.6 | -3 | 1046 | 1046 |
| 590 | 135 | 4.4 | 4.64 | -.30 | .07 | 0.9 | -1 | 0.8 | -1 | 1047 | 1047 |
| 514 | 118 | 4.4 | 4.34 | -.10 | .08 | 1.2 | 1 | 1.2 | 1 | 1048 | 1048 |
| 677 | 202 | 3.4 | 3.18 | .52 | .05 | 0.8 | -2 | 0.8 | -2 | 1049 | 1049 |
| 442 | 119 | 3.7 | 3.28 | .47 | .07 | 0.8 | -1 | 0.8 | -1 | 1050 | 1050 |
| 456 | 119 | 3.8 | 3.55 | .34 | .07 | 0.8 | -1 | 0.8 | -1 | 1051 | 1051 |
| 382 | 102 | 3.7 | 3.98 | .11 | .07 | 1.2 | 1 | 1.2 | 1 | 1052 | 1052 |
| 473 | 102 | 4.6 | 4.40 | -.13 | .09 | 0.9 | 0 | 0.8 | -1 | 1053 | 1053 |
| 545 | 152 | 3.6 | 3.91 | .15 | .06 | 0.9 | 0 | 0.9 | 0 | 1054 | 1054 |
| 578 | 170 | 3.4 | 3.74 | .24 | .06 | 1.0 | 0 | 0.9 | 0 | 1056 | 1056 |
| 469 | 119 | 3.9 | 4.56 | -.24 | .08 | 1.1 | 1 | 1.1 | 0 | 1057 | 1057 |
| 689 | 151 | 4.6 | 4.81 | -.44 | .07 | 1.1 | 0 | 0.9 | 0 | 1058 | 1058 |
| 444 | 100 | 4.4 | 4.37 | -.11 | .08 | 1.2 | 1 | 1.2 | 1 | 1059 | 1059 |
| 584 | 135 | 4.3 | 3.95 | .13 | .07 | 1.0 | 0 | 0.9 | 0 | 1060 | 1060 |
| 466 | 136 | 3.4 | 2.93 | .64 | .06 | 0.8 | -1 | 0.8 | -1 | 1061 | 1061 |
| 473 | 115 | 4.1 | 4.30 | -.07 | .07 | 1.3 | 2 | 1.3 | 2 | 1062 | 1062 |
| 366 | 85 | 4.3 | 3.84 | .19 | .09 | 0.9 | 0 | 0.9 | 0 | 1063 | 1063 |
| 319 | 67 | 4.8 | 4.62 | -.28 | .11 | 1.6 | 2 | 1.6 | 2 | 1064 | 1064 |
| **494** | **101** | **4.9** | **4.89** | **-.50** | **.10** | **2.0** | **4** | **2.1** | **4** | **1066** | **1066** |
| 575 | 153 | 3.8 | 3.65 | .28 | .06 | 1.2 | 2 | 1.2 | 2 | 1067 | 1067 |

Table 4 (continued)

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | ZStd | Outfit MnSq | ZStd | Num readers | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 358 | 85 | 4.2 | 3.89 | .16 | .09 | 1.2 | 1 | 1.2 | 1 | 1068 | 1068 |
| 608 | 136 | 4.5 | 4.71 | -.35 | .07 | 1.6 | 3 | 1.5 | 3 | 1069 | 1069 |
| 533 | 135 | 3.9 | 3.90 | .16 | .07 | 1.2 | 1 | 1.2 | 1 | 1070 | 1070 |
| 390 | 101 | 3.9 | 3.51 | .35 | .07 | 1.0 | 0 | 1.0 | 0 | 1071 | 1071 |
| 418 | 101 | 4.1 | 3.99 | .11 | .08 | 1.2 | 1 | 1.1 | 0 | 1072 | 1072 |
| 579 | 153 | 3.8 | 3.99 | .11 | .07 | 1.1 | 0 | 1.0 | 0 | 1073 | 1073 |
| 181 | 51 | 3.5 | 3.33 | .44 | .10 | 0.8 | -1 | 0.8 | -1 | 1074 | 1074 |
| 703 | 170 | 4.1 | 4.21 | -.02 | .06 | 0.9 | 0 | 0.8 | -1 | 1075 | 1075 |
| 579 | 118 | 4.9 | 4.63 | -.29 | .09 | 1.2 | 1 | 1.0 | 0 | 1076 | 1076 |
| 77 | 17 | 4.5 | 4.39 | -.13 | .21 | 0.6 | -1 | 0.6 | -1 | 1077 | 1077 |
| 413 | 102 | 4.0 | 3.94 | .14 | .08 | 0.7 | -2 | 0.7 | -2 | 1078 | 1078 |
| 475.3 | 115.3 | 4.1 | 4.14 | .00 | .08 | 1.0 | -0.1 | 1.0 | -0.4 | Mean | |
| 147.0 | 35.2 | 0.5 | 0.47 | .28 | .03 | 0.3 | 2.3 | 0.3 | 2.2 | S.D. | |

**RMSE (Model)** .09 **Adj S.D.** .27 **Separation** 3.10 **Reliability** .91
**Fixed (all same) chi-square: 874.6 d.f.: 65 significance: .00**
**Random (normal) chi-square: 65.5 d.f.: 64 significance: .43**

Table 5

*Summary Statistics for Items (arranged by the number of items)*

| Obsvd Score | Obsvd Count | Obsvd Average | Fair-M Avrage | Measure | Model S.E. | Infit MnSq | Infit ZStd | Outfit MnSq | Outfit ZStd | Nu items |
|---|---|---|---|---|---|---|---|---|---|---|
| 1842 | 450 | 4.1 | 4.13 | .03 | .04 | 0.9 | -1 | 0.9 | 0 | 1 1 |
| 1953 | 450 | 4.3 | 4.40 | -.13 | .04 | 0.7 | -5 | 0.7 | -4 | 2 2 |
| 2047 | 450 | 4.5 | 4.62 | -.28 | .04 | 0.7 | -5 | 0.7 | -4 | 3 3 |
| 1991 | 450 | 4.4 | 4.49 | -.19 | .04 | 0.7 | -5 | 0.7 | -5 | 4 4 |
| 1947 | 449 | 4.3 | 4.39 | -.13 | .04 | 0.8 | -3 | 0.9 | -2 | 5 5 |
| 2046 | 449 | 4.6 | 4.63 | -.29 | .04 | 0.6 | -5 | 0.7 | -5 | 6 6 |
| 1995 | 450 | 4.4 | 4.50 | -.20 | .04 | 0.6 | -6 | 0.7 | -5 | 7 7 |
| 1792 | 447 | 4.0 | 4.03 | .08 | .04 | 1.0 | 0 | 1.0 | 0 | 8 8 |
| 1896 | 444 | 4.3 | 4.33 | -.09 | .04 | 0.7 | -5 | 0.7 | -5 | 9 9 |
| 1768 | 440 | 4.0 | 4.05 | .08 | .04 | 0.7 | -5 | 0.7 | -4 | 10 10 |
| 1675 | 448 | 3.7 | 3.73 | .25 | .04 | 0.9 | -2 | 0.9 | -1 | 11 11 |
| 1892 | 448 | 4.2 | 4.27 | -.05 | .04 | 0.8 | -3 | 0.8 | -2 | 12 12 |
| 1917 | 448 | 4.3 | 4.33 | -.09 | .04 | 0.8 | -2 | 0.9 | -2 | 13 13 |
| 2015 | 445 | 4.5 | 4.60 | -.27 | .04 | 0.9 | -1 | 1.0 | 0 | 14 14 |
| **1581** | **447** | **3.5** | **3.49** | **.36** | **.04** | **2.0** | **9** | **1.9** | **9** | **15 15** |
| 1722 | 448 | 3.8 | 3.84 | .19 | .04 | 1.7 | 9 | 1.7 | 8 | 16 16 |
| **1289** | **448** | **2.9** | **2.73** | **.74** | **.04** | **2.1** | **9** | **2.0** | **9** | **17 17** |
| 1845.2 | 447.7 | 4.1 | 4.15 | .00 | .04 | 1.0 | -1.5 | 1.0 | -1.0 | Mean |
| 190.6 | 2.6 | 0.4 | 0.47 | .26 | .00 | 0.4 | 5.2 | 0.4 | 5.0 | S.D. |

**RMSE (Model)** .04 **Adj S.D.** .26 **Separation** 6.68 **Reliability** .98
**Fixed (all same) chi-square:** 819.2 **d.f.:** 16 **significance:** .00
**Random (normal) chi-square:** 16.0 **d.f.:** 15 **significance:** .38

## U.S. Department of Education
Office of Educational Research and Improvement (OERI)
National Library of Education (NLE)
Educational Resources Information Center (ERIC)

# REPRODUCTION RELEASE
(Specific Document)

ERIC
Educational Resources Information Center

## I. DOCUMENT IDENTIFICATION:

Title: A RASCH MEASUREMENT IN GRANT APPLICATION PROCESS

Author(s): Yesim CAPA & William E. LOADMAN

Corporate Source: The Ohio State University

Publication Date: APRIL 2003

## II. REPRODUCTION RELEASE:

In order to disseminate as widely as possible timely and significant materials of interest to the educational community, documents announced in the monthly abstract journal of the ERIC system, *Resources in Education* (RIE), are usually made available to users in microfiche, reproduced paper copy, and electronic media, and sold through the ERIC Document Reproduction Service (EDRS). Credit is given to the source of each document, and, if reproduction release is granted, one of the following notices is affixed to the document.

If permission is granted to reproduce and disseminate the identified document, please CHECK ONE of the following three options and sign at the bottom of the page.

The sample sticker shown below will be affixed to all Level 1 documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL HAS BEEN GRANTED BY

Sample

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

1

Level 1

Check here for Level 1 release, permitting reproduction and dissemination in microfiche or other ERIC archival media (e.g., electronic) *and* peper copy.

The sample sticker shown below will be affixed to all Level 2A documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE, AND IN ELECTRONIC MEDIA FOR ERIC COLLECTION SUBSCRIBERS ONLY, HAS BEEN GRANTED BY

Sample

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

2A

Level 2A

[X]

Check here for Level 2A release, permitting reproduction and dissemination in microfiche and in electronic media for ERIC archival collection subscribers only

The sample sticker shown below will be affixed to all Level 2B documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE ONLY HAS BEEN GRANTED BY

Sample

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

2B

Level 2B

Check here for Level 2B release, permitting reproduction and dissemination in microfiche only

Documents will be processed as indicated provided reproduction quality permits.
If permission to reproduce is granted, but no box is checked, documents will be processed et Level 1.

I hereby grant to the Educational Resources Information Center (ERIC) nonexclusive permission to reproduce and disseminate this document as indicated above. Reproduction from the ERIC microfiche or electronic media by persons other than ERIC employees and its system contractors requires permission from the copyright holder. Exception is made for non-profit reproduction by libraries and other service agencies to satisfy information needs of educators in response to discrete inquiries.

Sign here, → please

Signeture: Yesim Capa

Printed Neme/Position/Title: YESIM CAPA

Organization/Address: The Ohio State University 301 Ramseyer Hall. 29 W. Woodruff Ave Columbus, OH- 43210

Telephone: 688-5413

FAX:

E-Mail Address: capa.1@osu.edu

Date: 04/24/2003

*(Over)*

# III. DOCUMENT AVAILABILITY INFORMATION (FROM NON-ERIC SOURCE):

If permission to reproduce is not granted to ERIC, or, if you wish ERIC to cite the availability of the document from another source, please provide the following information regarding the availability of the document. (ERIC will not announce a document unless it is publicly available, and a dependable source can be specified. Contributors should also be aware that ERIC selection criteria are significantly more stringent for documents that cannot be made available through EDRS.)

| Publisher/Distributor: |
|---|
| Address: |
| Price: |

# IV. REFERRAL OF ERIC TO COPYRIGHT/REPRODUCTION RIGHTS HOLDER:

If the right to grant this reproduction release is held by someone other than the addressee, please provide the appropriate name and address:

| Name: |
|---|
| Address: |

# V. WHERE TO SEND THIS FORM:

Send this form to the following ERIC Clearinghouse:   University of Maryland
ERIC Clearinghouse on Assessment and Evaluation
1129 Shriver Lab, Bldg 075
College Park, MD 20742
Attn: Acquisitions

However, if solicited by the ERIC Facility, or if making an unsolicited contribution to ERIC, return this form (and the document being contributed) to:

University of Maryland
ERIC Clearinghouse on Assessment and Evaluation
1129 Shriver Lab, Bldg 075
College Park, MD 20742
Attn: Acquisitions